

RESEARCH ARTICLE

Advancing Oncology Diagnostics: AI-Enabled Early Detection of Lung Cancer Through Hybrid Histological Image Analysis

NAGLAA F. NOAMAN¹, BASSAM M. KANBER¹, AHMAD AL SMADI²,
LICHENG JIAO¹, (Fellow, IEEE), AND MUTASEM K. ALSMADI³, (Member, IEEE)

¹School of Artificial Intelligence, Xidian University, Xian 710071, China

²Department of Data Science and Artificial Intelligence, Zarqa University, Zarqa 13100, Jordan

³Department of Management Information Systems, College of Applied Studies and Community Service, Imam Abdulrahman Bin Faisal University, Dammam 34212, Saudi Arabia

Corresponding author: Licheng Jiao (lchjiao@mail.xidian.edu.cn)

This work was supported in part by the Key Scientific Technological Innovation Research Project by Ministry of Education; in part by the Joint Funds of the National Natural Science Foundation of China under Grant U22B2054; in part by the National Natural Science Foundation of China under Grant 62076192, Grant 61902298, Grant 61573267, Grant 61906150, and Grant 62276199; in part by the 111 Project; in part by the Program for Cheung Kong Scholars and Innovative Research Team in University under Grant IRT 15R53; in part by the ST Innovation Project from the Chinese Ministry of Education; in part by the Key Research and Development Program in Shaanxi Province of China under Grant 2019ZDLGY03-06; and in part by China Postdoctoral Fund under Grant 2022T150506.

ABSTRACT Against the backdrop of the pervasive global challenge of cancer, with particular emphasis on lung cancer (LC), this study centers its investigation on the critical realm of early detection leveraging artificial intelligence (AI) within the domain of histological image analysis. Through the fusion of DenseNet201 with color histogram techniques, a novel hybrid feature set emerges, engineered to elevate classification accuracy. The comprehensive evaluation encompasses eight diverse machine learning (ML) algorithms, spanning from K-Nearest Neighbors (KNN) to Support Vector Machines (SVM), including notable contenders such as LightGBM (LGBM), CatBoost, XGBoost, decision trees (DT), random forests (RF), and multinomial naive Bayes (MultinomialNB). This rigorous examination illuminates a distinguished model, achieving a remarkable accuracy rate of 99.683% on the LC25000 dataset. The extension of this methodology to breast cancer detection, utilizing the BreakHis dataset, yields a commendable accuracy rate of 94.808%. These findings underscore the transformative potential of AI in the intricate landscape of histopathological analysis, positioning it as a pivotal force in advancing diagnostic capabilities. A meticulous comparative analysis not only underscores the merits but also elucidates the limitations of existing AI applications in medical imaging, thereby charting a roadmap for future refinements and clinical deployments. Consequently, continued research in AI within clinical settings is advocated, with the ultimate aim of fortifying early cancer diagnosis and subsequently enhancing patient outcomes through judicious therapeutic interventions.

INDEX TERMS Densenet201, histopathological images, image processing, lung cancer, machine learning.

I. INTRODUCTION

In the year 2023, the United States is grappling with a heavy and sobering truth: an estimated 1,958,310 new cancer cases and 609,820 cancer-related deaths are anticipated, casting

The associate editor coordinating the review of this manuscript and approving it for publication was Yeliz Karaca.

a shadow of concern over our nation's health landscape. This translates to about 5,365 new cancer cases and 1,671 cancer-related deaths each day, each a poignant reminder of the ongoing struggle against this relentless disease [1]. Among these troubling statistics, lung cancer (LC) emerges as a formidable opponent, with the American Cancer Society predicting approximately 238,340 new cases in the United

States for the year 2023 [2]. However, LC is often diagnosed in our elderly population, with an average age of 70 [3]. This paints a poignant picture of the challenges our elderly population faces in the battle against this formidable enemy. LC is a major health concern for people all around the world. It is the top cause of cancer-related deaths in various communities, and its impact on healthcare, families, and economies is substantial. Even with advancements in medical technology, the number of LC cases is increasing globally [4], [5]. This emphasizes the critical importance of creating new and effective ways to detect and treat the disease on a global scale. The research's focus on using pretrained Densenet201, color histogram, and machine learning (ML) for diagnostics tackles an urgent health issue and plays a part in the larger effort to enhance cancer care.

The sheer magnitude of this predicament is underscored by the fact that LC has the distinction of being the leading cause of cancer-related deaths in the United States, claiming the lives of approximately one in five individuals who succumb to this devastating disease [6]. Annually, it causes losses greater than the cumulative impact of colon, breast, and prostate cancer combined, reinforcing its unparalleled impact on public health [2], [7]. The development of cancer is a complex interaction and combination of behavioral and environmental factors. Smoking, obesity, alcohol abuse, radiation exposure, and biological factors are all known factors in the development of cancer [8]. The challenges associated with early cancer detection are particularly stressful. It often remains asymptomatic or shows only subtle signs in its early stages, and is therefore difficult to detect. By the time symptoms become manifest, the tumor has usually reached an advanced stage, making timely diagnosis a formidable task [9].

Currently, when it comes to LC detection, it mainly uses imaging methods like chest X-rays and CT scans, along with tissue biopsies for detailed analysis [10]. While these techniques have been really helpful in diagnosing LC, they do have significant limitations. Imaging techniques can miss early-stage tumors or non-solid nodules, leading to false negatives. Moreover, how these tests are interpreted relies a lot on the skill of the radiologists and pathologists, which can introduce some subjectivity and variation in the diagnosis. The significance of this research lies in its ability to change LC detection methodologies through the use of pretrained Densenet201, color histogram, and ML techniques to analyze detailed images of lung tissue. The goal is to enhance the accuracy, speed, and efficiency of diagnosing LC in its early stages. They can automate the detection process, reducing the reliance on manual examination and potentially decreasing the rate of misdiagnosis. This kind of automation is super important for diagnosing LC before it gets really serious, which can make a big difference in how well the treatment works and the chances of surviving. Building on the methodologies and findings of this research, future studies could explore the adaptation of our model for other types of cancer, leveraging the unique patterns and characteristics

present in different cancerous tissues. Also integrating these models into clinical settings, where they could be tested and refined in real-world scenarios. This integration would validate the efficacy of the models in practical applications, streamline diagnostic processes, improve patient outcomes through earlier detection, and tailor treatment plans more effectively.

In this sobering landscape, the need for innovative solutions has become increasingly urgent. The integration of diverse feature extraction (FE) techniques with pretrained models and ML to improve the accuracy and reliability of LC diagnosis from histopathological images has been a focal point of extensive research [11], [12], [13]. Our research sets out to harness the boundless potential of AI in the domain of histological image analysis, with a specific focus on LC detection. The journey begins with a recognition of the basic dynamics of cancer development—where some cells, either damaged or reaching the end of their lifespan, fail to be replaced by healthy cells, resulting in the formation of tumors [14]. These tumors may take the form of benign or malignant growths, with the latter characterized by aggressive, abnormal cell proliferation that can rapidly invade and damage surrounding tissues [15]. This research comes to the fore armed with the LC25000 dataset. It encompasses histological images across five distinct classes, including colon adenocarcinomas, benign colonic tissues, lung adenocarcinomas, lung squamous cell carcinomas, and benign lung tissues. Our mission is to analyze these images by using advanced AI methodologies, with the goal of saving lives by identifying LC in its earliest, most treatable stages.

A. RELATED WORKS

Digital pathology has revolutionized cancer diagnosis, especially lung cancer (LC) classification using histopathological images. The wealth of research on image analysis, from enhancement and feature extraction (FE) to deep learning (DL) and machine learning (ML) methods for reliable classification, reflects this development. The following sections discuss how many key research techniques, datasets, and results have shaped the classification of LC histopathology images. Therefore, based on existing knowledge, a complete evaluation of these linked studies will contextualize our findings within the larger scientific debate, recognizing both advances and limitations.

Al-Jabbar et al. [11] proposed multiple strategies for histological image classification of lung and colon cancer. One of them involved image enhancement through filtering and contrast improvement, FE using DL models (GoogLeNet and VGG-19), and the integration of these features with handcrafted methods. This particular approach demonstrated high accuracy, with the fusion features of VGG-19 and handcrafted methods achieving the highest accuracy rate of 99.64%. Garg and Garg [12] employed eight distinct pre-trained CNN models, including VGG16, NASNetMobile, InceptionV3,

InceptionResNetV2, ResNet50, Xception, MobileNet, and DenseNet169, to classify cancerous and non-cancerous images from the LC25000 dataset. The models consistently achieved notable results with accuracies ranging from 96% to 100%. The use of data augmentation and attention visualization techniques, such as GradCAM and SmoothGrad, contributed to the understanding of models' decision processes and strengthens the potential for automated cancer diagnosis. Kumar et al. [13] contrasted six conventional FE techniques (fuzzy color and texture histogram, color correlogram, color layout, edge histogram, pyramid histogram of oriented gradients, pyramid based local binary patterns) with transfer learning approach, utilizing seven pre-trained CNN models for FE from lung and colon histopathological images. The transfer learning approach, particularly with DenseNet-121, outperformed conventional classifiers, achieving 98.60% accuracy, 98.63% precision, 98.60% recall, 0.985 f1-score, and 0.1 ROC-AUC. The methodology involved image resizing and pre-processing, emphasizing the importance of color and texture-based features in cancer detection.

Masud et al. [16] proposed a classification framework for LC. They utilized unsharp masking (UM) to enhance contrast at color junctions. This technique sharpens the original image by subtracting a blurred version from itself. FE was performed using the extraction of 2D fourier features and the extraction of 2D wavelet features. The extracted features were then processed and fed into a convolutional neural network (CNN) for classification. Utilizing the LC25000 dataset for this purpose, the results demonstrated remarkable accuracy of up to 96.33%. Ali and Ali [17] proposed a novel multi-input dual-stream capsule network employing two convolutional layer blocks: the convolutional layers block (CLB), which uses traditional convolutional layers, and the separable convolutional layers block (SCLB), which employs separable convolutional layers. The CLB processes unprocessed histopathology images, while the SCLB handles uniquely pre-processed images using techniques like color balancing, gamma correction, image sharpening, and multi-scale fusion. The empirical analysis on the LC25000 dataset showed significant improvement in classification results, with achieving 99.58% overall accuracy. Baranwal et al. [18] utilized DL and CNNs for LC diagnosis, employing the LC25000 dataset with 15,000 images. Four CNN models—ResNet50, VGG-19, Inception-ResNet-V2, and DenseNet121—were used for classifying LC types. The study emphasized the effectiveness of DenseNet121, especially when used with triplet neural networks, achieving an accuracy of 99.08% in differentiating between cancerous and non-cancerous tissues.

Hatuwal and Thapa [19] leveraged CNNs to classify histopathological images of lung tissue, including benign tissue, adenocarcinoma, and squamous cell carcinoma. The study encompassed various stages, from using the LC25000 histopathological image dataset, data formatting for uniformity and efficiency, to model training, testing, and prediction.

The CNN architecture is composed of convolutional layers, max pooling, and fully connected layers, with a focus on optimizing image classification. Results indicated an impressive training accuracy of 96.11% and validation accuracy of 97.20%, demonstrating the potential of CNNs in enhancing LC diagnosis. Mehmood et al. [20] developed an effective model for diagnosing lung and colon cancers using the LC25000 dataset of histopathology images. They initially achieved 89% accuracy by employing a modified AlexNet neural network. To further improve performance, especially for an underperforming class, they applied a contrast enhancement technique, which boosted the overall accuracy to 98.4%. Mangal et al. [21] focused on creating a computer-aided diagnosis system using CNNs to identify lung and colon cancer from digital pathology images. They utilized the LC25000 dataset to classify histopathological slides into various cancer types and benign tissues. The study demonstrated the potential of CNNs in achieving high diagnostic accuracy, with rates exceeding 97% for LC and over 96% for colon cancer.

Civit-Masot et al. [22] focused on designing, implementing, and evaluating a diagnostic aid system for non-small cell LC (NSCLC) detection using DL. The study included an explainable DL component that informs pathologists about the image areas used for classification and the confidence of each class utilizing LC25000 dataset. The system showed high accuracy between 97.11% and can potentially reduce the time spent by pathologists on each patient, thereby speeding up diagnostic processes. Mamun et al. [23] focused on developing a model to predict LC using ensemble learning methods, namely XGBoost, LightGBM, Bagging, and AdaBoost. They evaluated these techniques on a dataset of 309 individuals collected from kaggle, considering factors like age, smoking habits, and symptoms such as fatigue, allergy, and chest pain. The study found XGBoost to be the most effective, with an accuracy of 94.42%.

Ramesh et al. [24] developed a multi-level CNN (ML-CNN) architecture for detecting different types of LC. They employed a multi-scale convolution strategy to effectively extract features from lung nodules of various sizes and morphologies. The model was evaluated using the LC25000 dataset, which includes histopathological images of squamous cell cancer and adenocarcinoma. The model demonstrated superior performance compared to traditional methods, achieving an accuracy of 64% in training and 89% in validation. Shanmugam and Rajaguru [25] presented a novel methodology for detecting LC using histopathological images. Their approach focused on preprocessing and segmentation, followed by FE using particle swarm optimization (PSO) and grey wolf optimization (GWO). They also utilized algorithms such as KL divergence and invasive weed optimization (IWO) for feature selection. The study employed seven classifiers to classify the images into benign or malignant, achieving an impressive accuracy of 91.57% with the DT classifier. This high accuracy

was attained by using GWO for FE, IWO for feature selection, and the RAdam approach for hyperparameter tuning. Krishnan et al. [26] introduced an improved graph neural network (IGNN) optimized by the green anaconda optimization (GAO) algorithm to maximize accuracy in segmenting and classifying LC. The process involved pre-processing images using the gabor filter method, segmentation with the modified expectation maximization (MEM) algorithm, and FE through the histogram of oriented gradient (HOG) scheme.

These works demonstrate how the field of LC histological image classification is changing, emphasizing the efficacy of combining DL, conventional, and hybrid approaches. The relevant literature has been carefully reviewed, including an assessment of its efficacy, strengths, and drawbacks. The findings have been consolidated and displayed in Table 1. Our study uses many ML algorithms, including KNN, LGBM, CatBoost, XGBoost, DT, RF, MultinomialNB, and SVM. To the best of our knowledge, this study is the initial endeavor to compare these techniques for LC histopathology images, covering the analysis of binary classes. In addition, the research combines DenseNet201 and color histogram, establishing a strong basis for the investigation.

B. CONTRIBUTIONS

The proposed work extends beyond traditional methods by exploring a fusion of features extracted through a combination of approaches. The high-dimensional and nuanced data present in these histological images demand a novel perspective. Early-stage abnormal cells often share similar characteristics, leading to the development of hybrid systems that amalgamate features from diverse sources. This enhances the ability to discriminate between subtle variations that may signify the onset of LC [27]. The following paragraphs provide an overview of the contributions made in the study:

- DenseNet201 and color histogram methods are utilized to extract informative features from histological images within the LC25000 dataset and combine them, with a primary focus on lung adenocarcinomas, lung squamous cell carcinomas, and benign lung tissues. This results in a hybrid feature set that enhances the discriminative power of the classification models.
- The performance of eight machine learning (ML) algorithms is assessed, including K-Nearest neighbors (KNN), LGBM, CatBoost, XGBoost, decision trees (DT), random forests (RF), multinomial naive bayes (MultinomialNB), and support vector machines (SVM). The most effective algorithm is identified and selected through a comprehensive evaluation encompassing accuracy, specificity, precision, recall, f1-score, and computational efficiency.
- Multi-class classification is conducted to differentiate between lung adenocarcinomas, lung squamous cell carcinomas, and benign LC-related classes. Subsequently, binary classification tasks are engaged, including benign

versus adenocarcinomas, benign versus carcinomas, and adenocarcinomas versus carcinomas, to deepen the understanding of LC subtypes and their distinctive characteristics.

The remainder of the paper is organized into the following sections: Methodology, results and analysis, discussion, and conclusion. These sections collectively present a perspective on a future scenario in which LC diagnostics are not only precise but also readily accessible, timely, and revolutionary.

II. METHODOLOGY

This section delineates the methodology for lung cancer (LC) classification, encapsulating critical stages integral to the approach. The framework comprises sequential phases: dataset utilization, data preprocessing, feature extraction (FE), feature combination, application of machine learning (ML) algorithms, and evaluation metrics. Each phase plays a pivotal role in the overarching classification process, contributing uniquely to enhancing accuracy and efficiency in LC detection.

A. LC25000 DATASET

For the purposes of our research on LC, we focused on the LC25000 dataset, a comprehensive collection of histopathological images. This dataset, originating from the Kaggle platform, was compiled by Andrew Borkowski and his team at James Hospital in Tampa, Florida. It is segmented into various cancer types, including both colon and (LCs). Out of the total 25,000 images in the dataset, we selectively utilized 15,000 images representing three LC categories: adenocarcinoma (lung_aca), which forms a considerable portion of LC cases, benign lung tissue (lung_bnt), and squamous cell carcinoma (lung_scc), the second most common type, ensuring a targeted approach to the proposed study, and allowing to delve deeper into the specificities of LC. Each category comprises an equal number of images, 5,000 per type, thus maintaining a balance in our analysis. Subsequently, Fig. 1 shows samples of the dataset, and a pie chart is employed to visualize the distribution of the various classes in our dataset in Table 2. Visualization is crucial for understanding the dataset's composition, ensuring that our model is trained and tested on a balanced and representative sample. Originally, this dataset was derived from 1,250 primary images (250 for each type) collected from pathology slides. These were augmented using various techniques like rotation and flipping, expanding the total to 25,000 images. These images, each meticulously prepared and cropped from their original dimensions of 1024×768 pixels to a uniform size of 768×768 pixels [28].

B. PREPROCESSING STAGE

The preprocessing of the LC image dataset is a critical step in our study. The initial step involves loading images with the RGB color mode. These images are scaled to 128×128 pixels to balance computational efficiency and

TABLE 1. Comprehensive comparative analysis of efficacy, strengths, and drawbacks among diverse imaging techniques applied to the diagnosis and treatment of Lung cancer (LC).

Ref	Year	Dataset	Efficacy	Strengths	Drawbacks	Hardware
[11]	2023	LC25000	VGG-19 + Handcrafted features + ANN 99.64% accuracy, 99.85% sensitivity, 100% specificity and precision.	<ul style="list-style-type: none"> Hybrid AI systems integrating CNN models with handcrafted features. VGG-19 + Handcrafted was optimal. Applied on colon and lung. High performance metrics, indicating the model's reliability. 	<ul style="list-style-type: none"> Single dataset, results may not generalize to other datasets. Time computation not available. Huge number of features. 	Not provided
[12]	2020	LC25000	Pre-trained CNN models with visualization of class activation and saliency maps accuracy of 96-100% in classifying malignant vs benign tumors.	<ul style="list-style-type: none"> Application of visualization techniques (GradCAM, SmoothGrad) for interpretability. Effective use of pre-trained CNN models. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. Binary classification. Lack of detailed about computational environment. 	Not provided
[13]	2022	LC25000	DenseNet121 FE + RF 98.60% accuracy, 98.63% precision, 98.60% recall, an f1-score of 0.986.	<ul style="list-style-type: none"> Evaluates classifier performance using multiple metrics. Applied on colon and lung. Comprehensive comparison FE. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. 	Python 3.8, IBM Intel Core i-7-6700 CPU @ 3.40 GHz processor, 8 GB RAM, NVIDIA GeForce GPU.
[16]	2021	LC25000	Unsharp masking for image sharpening: 2D Fourier and wavelet transforms for FE + CNN Model 96.33% accuracy, 96.39 % precision.	<ul style="list-style-type: none"> Uses fourier and wavelet transforms to extract complementary feature sets. Enhance CNN by Employing a custom-designed 4-channel CNN architecture. Applied on colon and lung. High performance metrics. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. Computational resources. 	Not provided
[17]	2021	LC25000	Capsule network + conventional and separable CNNs 99.58% accuracy, 98.66% precision.	<ul style="list-style-type: none"> Allows the model to learn features from both unprocessed and pre-processed images. Use of capsule networks with convolutional layers. Improve the overall feature learning process of the model. Applied on colon and lung. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. Computational complexity due to the dual-input approach. 	Windows 10 PC, Nvidia GeForce GTX 1060, 16 GB RAM, Intel C17 64-bit, Keras and TensorFlow.
[18]	2021	LC25000	CNN model using triplet loss 99.08% accuracy with DenseNet121.	<ul style="list-style-type: none"> Comprehensive exploration of various CNN architectures. Application of triplet loss improves the differentiation between the classes. Applied on colon and lung. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. Specific hardware details are not provided affect reproducibility. 	Python
[19]	2020	LC25000	CNN model 96.11% training accuracy, 97.20% validation accuracy.	<ul style="list-style-type: none"> Improving the quality of input data for the CNN model by image pre-processing. 	<ul style="list-style-type: none"> Single dataset. Lack comparison with other algorithms. Specific hardware details are not provided affect reproducibility. 	Google Colaboratory GPU
[20]	2022	LC25000	Initial accuracy 89%, improved to 98.4% by AlexNet CNN +CISP (Histogram Equalization).	<ul style="list-style-type: none"> Improves classification using CSIP. Minimizes computational costs. Applied on colon and lung. Computational efficiency. 	<ul style="list-style-type: none"> Single pretrained model. Single dataset. Limited description of CSIP. Time computation not available. 	Not provided
[21]	2020	LC25000	Shallow CNN 97.92% and 96.95% accuracy (lung and colon respectively).	<ul style="list-style-type: none"> Integration of DL. Develop CNN models. Effective use of shallow CNN architecture. 	<ul style="list-style-type: none"> Single dataset. Lacks some implementation details. Lack of detailed about computational environment. 	Google's Colab TensorFlow
[22]	2022	LC25000	Train CNN model and used explainable DL (GradCAM) 97.11% accuracy.	<ul style="list-style-type: none"> Highly accurate, explainable. Applied explainable DL techniques. Highlighting specific image areas used for classification. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. 	Google's Colab TensorFlow
[23]	2022	Kaggle	Ensemble learning techniques (XGBoost, LightGBM, bagging, and AdaBoost) 94.42% accuracy.	<ul style="list-style-type: none"> Effectiveness of XGBoost in LC prediction. Comprehensive evaluation of ensemble learning techniques. 	<ul style="list-style-type: none"> Single dataset and small. Time computation not available. 	Not provided
[24]	2023	LC25000	Multi-level CNN (ML-CNN) training accuracy: 64%, validation accuracy: 89%.	<ul style="list-style-type: none"> Handle the heterogeneity in lung nodule sizes and morphologies. Leveraging multi-scale convolution for improved FE. 	<ul style="list-style-type: none"> Single dataset and small. Time computation not available. 	Python 3.X and Google Colab. provided a Jupyter notebook - GPU
[25]	2023	LC25000	Grey wolf optimization (GWO) + Invasive Weed optimization (IWO) + hyperparameter tuning RAdam + DT accuracy of 91.57%.	<ul style="list-style-type: none"> Integration of PSO and GWO for FE. The use of hyperparameter tuning methods to improve accuracy. 	<ul style="list-style-type: none"> Single dataset. Time computation not available. Binary classification. Lack of detailed about computational environment. 	Not provided
[26]	2023	LC25000	Histogram of oriented gradient (HOG) + hyperparameter tuning green anaconda optimization (GAO) + improved graph neural network (IGNN) accuracy of 98.9%.	<ul style="list-style-type: none"> Employed gabor filter for pre-processing and MEM for segmentation. Introduced a novel IGNN model optimized by GAO. 	<ul style="list-style-type: none"> Single dataset. Lack of detailed about computational environment. 	Not provided

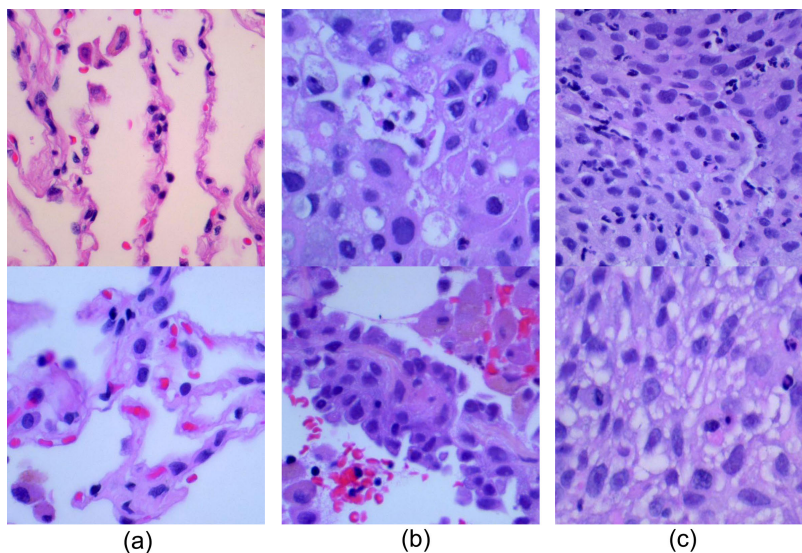


FIGURE 1. Sample of the LC25000 dataset images. (a) Benign lung tissues samples, (b) Lung adenocarcinomas samples, (c) Lung squamous cell carcinomas.

TABLE 2. Distribution of the three classes (Lung Adenocarcinomas, Lung Squamous cell Carcinomas and Benign Lung tissues) in LC25000 dataset.

Class Name	Number of Images	Percentage
Lung_aca	5000	33.3%
Lung_scc	5000	33.3%
Lung_n_benign	5000	33.3%

maintain sufficient detail for classification accuracy. Scaled images and corresponding labels are transformed into numpy arrays for further processing. Additionally, a proprietary function (ensure_correct_depth) was used, to standardize image depth to CV_8U. This step normalizes image data, ensuring that all images have the same scale and format for uniform FE throughout the dataset. Each image is converted to the HSV color space, for analyzing hue-based color distributions. Using the preprocess_input function, designed to preprocess images in the same manner as the original model training, we perform normalization of pixel values. This typically involves subtracting the mean RGB values and dividing by the standard deviation, specific to the ImageNet dataset. This normalization aligns the distribution of the input data with that used during the model’s training.

C. FEATURE EXTRACTION (FE)

The model presented in the suggested methodology is based on many crucial phases for the classification of LC, as seen in Fig 2. The approach employed herein hinges upon FE. The employed methodology utilizes a fusion of image processing and DL methodologies to extract important data from the histopathological images in the LC25000 dataset:

1) COLOR HISTOGRAM ANALYSIS

After converting images to the HSV (Hue, Saturation, Value) color space, which is preferred in image processing tasks as it separates image intensity (Value) from color information (Hue and Saturation) making it more resilient to changes in lighting conditions [29]. A color histogram is computed for each image, focusing on the Hue component, which represents the color type in the image and is important for capturing color-based features in medical images. This is needed to understand the color distribution within the images, which can be indicative of various lung conditions. The computed histograms are normalized to ensure uniformity in feature scaling, which is important as it brings all features to a comparable range, thereby preventing features with larger numeric ranges from dominating the learning process in the classification models [30]. This process converts the histogram into a probability distribution of Hue values. The normalized histogram of each image is flattened into a one-dimensional feature vector, making it suitable for use in ML models.

2) CONTOUR FE

After converting images from RGB to grayscale, which simplifies the image data and focuses on the structural information [31]. Contours in the grayscale images are detected. These objects correspond to areas of potential medical interest in lung scans. This is needed to focus on the significant shapes in the image, which are likely to be of medical relevance. The contour approximation is used to compress horizontal, vertical, and diagonal segments in the contour, thereby reducing the number of points required to represent a contour [32]. For each contour, various features are extracted, including the area and the perimeter of the contour, which are flattened to form a feature vector for each

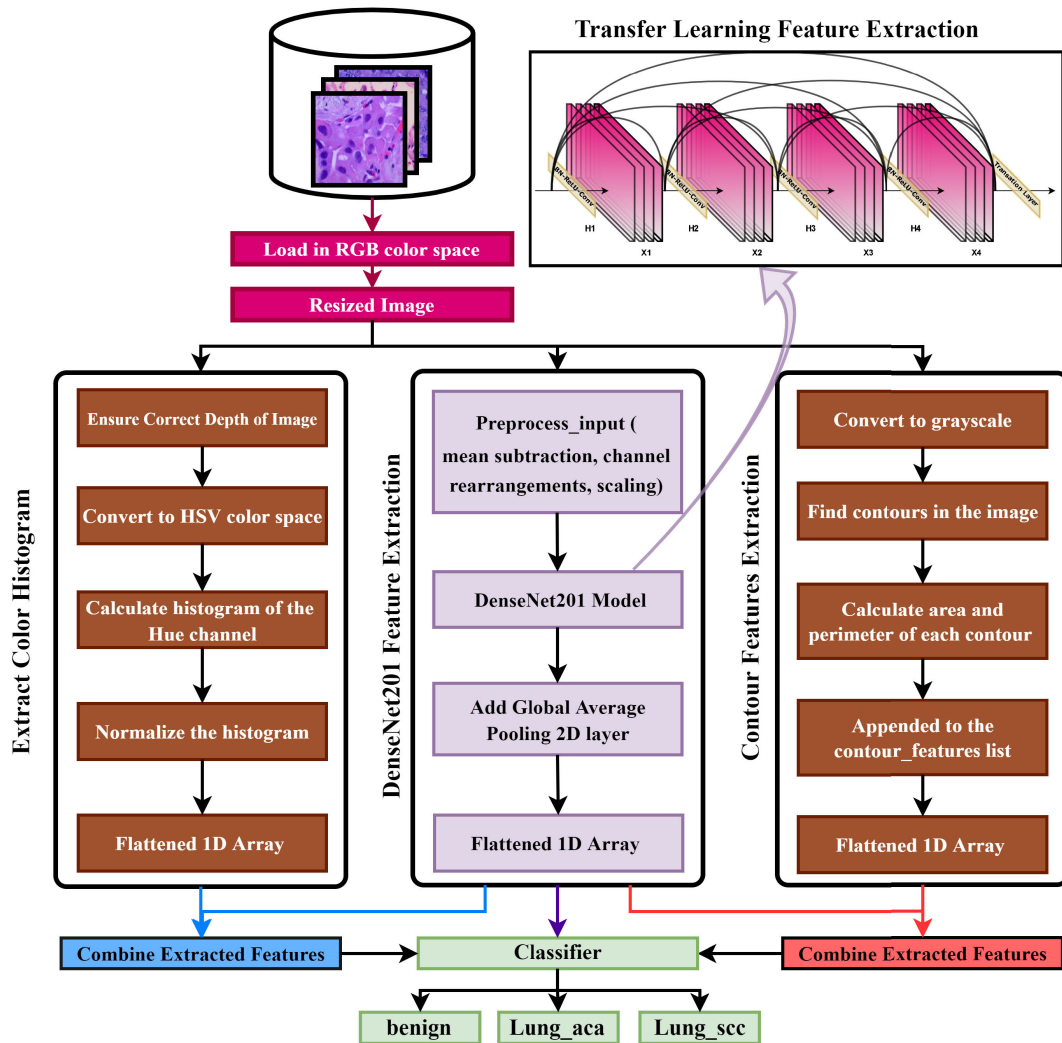


FIGURE 2. Structure of the model.

image. These flattened feature vectors represent the structural characteristics of the lung images, which are important for the subsequent classification process.

3) DENSENET201 FEATURES

DenseNet201 is a form of the Dense Convolutional Network architecture, which connects each layer to every other layer in a feed-forward fashion, to ensure maximum information flow between layers in the network [33]. We utilized it with pre-trained weights from the ImageNet database, to allow the model leverages knowledge from a vast dataset, enhancing its FE capabilities for the task of LC image classification. It is modified for FE by removing the top classification layer, which transforms it into a powerful feature extractor. As the images propagate through the network, the DenseNet201 layers extract a wide range of features, from basic edge features in the initial layers to more complex patterns in the deeper layers [34]. The final layer produces a

comprehensive feature representation of each image. After the last convolutional layer, a Global Average Pooling (GAP) layer reduces feature map dimensionality while keeping the most important spatial information. Then we result in a single feature vector for each image, encapsulating the important features identified by the network, which represent a deep and nuanced understanding of the image content, capturing both local and global patterns relevant to LC classification.

D. FEATURE COMBINATION

In our methodology, we employ specific strategies to combine the extracted features. Our focus is on leveraging the strengths of each feature type—DenseNet201, contour, and color histogram features—by creating combined feature sets for classification. Initially, we explore results of DenseNet201 as base learner. The second strategy, we utilize features extracted solely by the DenseNet201 model with ML classifiers. The third strategy, we explore combinations of

DenseNet201 features with contour features. This combination aims to harness both the high-level pattern recognition capabilities of DenseNet201 and the geometric information provided by contour analysis. The fourth strategy, we explore combinations of DenseNet201 features with color histogram features. This set integrates DL patterns with color textural information. By exploring different combinations, we can assess which feature sets are most effective for classifying LC images.

E. MACHINE LEARNING (ML) ALGORITHMS

In the domain of LC image classification, selecting appropriate ML models is pivotal for achieving high accuracy and reliability. Our methodology incorporates a diverse array of classifiers, each with unique strengths and characteristics. Below is a short description of each model selected for this study:

1) K-NEAREST NEIGHBORS (KNN)

KNN is a simple, effective instance-based learning method. It classifies samples through how similar they are to the training set. KNN is simple yet effective, especially when decision border is irregular [35].

2) LIGHT GRADIENT BOOSTING MACHINE (LGBM)

LGBM employs tree-based learning methods for gradient boosting. Highly efficient and performant, especially with huge datasets. LGBM excels at handling unbalanced data, which is prevalent in medical imaging datasets [36].

3) CATBOOST

CatBoost is another gradient-boosting technique optimized for categorical data. Its robustness and comprehensive feature combination handling make it ideal for our non-categorical data. Automated missing data management is useful in complicated datasets [37].

4) EXTREME GRADIENT BOOSTING (XGBOOST)

XGBoost is a quite efficient and scalable implementation of gradient boosting. Its performance in ML contests has made it popular. Speedy and powerful, XGBoost provides fine-grained model tweaking control [38].

5) DECISION TREES (DT)

DT is a simple and interpretable model that partitions data at each node depending on criteria. Their usefulness is in comprehending decision-making. It can help medical imaging professionals identify classification-relevant characteristics [39].

6) RANDOM FOREST (RF)

The ensemble learning technique RF builds many DTs during training. Overfitting, a problem with single DTs, is reduced, improving classification performance. This robust model handles linear and non-linear data well [40].

TABLE 3. Hyperparameters Configuration of the ML classifiers (α learning rate, l leaves, n estimators, d Max Depth, cbt colsample_bytree, k n_neighbors).

Classifier	Hyperparameters
LightGBM	$\alpha=0.09$, $n=200$, $l=31$, $d=-1$, $\text{min_child_samples}=20$, $\text{subsample}=1.0$, $\text{colsample_bytree}=1.0$
RF	$n=250$, $d=None$, $\text{min_samples_split}=2$, $\text{min_samples_leaf}=1$
DT	$d=None$, $\text{min_samples_split}=2$, $\text{min_samples_leaf}=1$
XGBoost	$n=300$, $d=6$, $\alpha=0.3$, $\text{subsample}=1$, $\text{colsample_bytree}=1$
CatBoost	$\text{verbose}=0$ $\text{iterations}=1000$, $\alpha=0.03$, $d=6$, $\text{l2_leaf_reg}=3$, $\text{loss_function}='Logloss'$
KNN	$k=4$, $\text{weights}='distance'$, $\text{metric}='manhattan'$, $\text{algorithm}='auto'$, $\text{leaf_size}=30$, $p=2$, $\text{metric_params}=None$
Naive Bayes	$\alpha=1.0$, $\text{fit_prior}=True$, $\text{class_prior}=None$
SVM	$\text{probability}=True$, $C=1.0$, $\text{kernel}='rbf'$, $\text{degree}=3$, $\text{gamma}='scale'$, $\text{coef0}=0.0$

7) MULTINOMIAL NAIVE BAYES (MULTINOMIALNB):

Naive Bayes' multinomialNB version handles multinomially distributed data. It assumes feature independence in image classification, simplifying computation. This approach works well for classification issues with huge feature sets, making it suited for our high-dimensional data [41].

8) SUPPORT VECTOR MACHINE (SVM)

SVM is a sophisticated classifier that finds the optimum hyperplane to divide classes in feature space. It's versatile and works well with high-dimensional data and linear and non-linear data [42].

The algorithms are applied to the feature vectors from the LC25000 dataset using the feature combination strategies mentioned above, and their performance in classifying images into lung adenocarcinomas, lung squamous cell carcinomas, and benign lung is critically evaluated.

F. HYPERPARAMETERS SELECTION

In the development of our predictive models, hyperparameter selection plays a crucial role in optimizing performance and achieving robustness. The choice of hyperparameters can significantly affect the learning process and the resulting model's efficacy. The research was conducted in a standardized computing setting to provide consistency and clarity. The simulations were conducted in a Google Colab environment (with 12 GB of RAM, 78 GB of HDD, and a cloud GPU). For each algorithm utilized, hyperparameters were meticulously selected to balance the trade-off between training time and model accuracy. Table 3 outlines the algorithms employed, along with the hyperparameters explicitly set for our experiments. This table demonstrates our methodical approach to hyperparameter selection, ensuring each model is finely tuned for optimal performance within the scope of our research objectives.

TABLE 4. The performance metrics results of the DenseNet201 base learner on the test set. MA (Micro Average), WA (Weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
DenseNet201	0.9733	0.4934	6908.51	0.97 MA	0.97 MA	0.97 MA
				0.97 WA	0.97 WA	0.97 WA

G. EVALUATION METRICS

The models will be evaluated based on various metrics, including accuracy, average specificity, time (S), precision, recall, and f1-score, as in (1) to (5). These metrics provide a comprehensive view of each model's performance, taking into account factors such as generalizability, efficiency, and the balance between precision and recall.

$$\text{Accuracy} = \frac{TP + TN}{TN + TP + FP + FN}. \quad (1)$$

$$\text{Precision} = \frac{TP}{TP + FP}. \quad (2)$$

$$\text{Recall} = \frac{TP}{TP + FN}. \quad (3)$$

$$F1 = 2 \times \frac{\text{Recall} \times \text{precision}}{\text{Recall} + \text{precision}}. \quad (4)$$

$$\text{Specificity} = \frac{TN}{TN + FP}. \quad (5)$$

III. RESULTS AND ANALYSIS

The study on the potential of ML algorithms in the multi-classification of LC using the LC25000 dataset revealed significant insights. The models were evaluated based on various metrics, including accuracy, specificity, precision, recall, and f1-score. The evaluation involved testing the models under three FE scenarios: using DenseNet201 alone, DenseNet201 combined with contour features, and DenseNet201 combined with histogram features.

A. PERFORMANCE OF DENSENET201 AS BASE LEARNER

Utilizing DenseNet201 as a standalone base learner yielded an impressive accuracy of 97.33% and an f1-score averaging 97%. Despite its high performance, the specificity stood at 49.34%, indicating room for improvement in distinguishing between classes. The model required 6908.51 seconds for execution, highlighting a trade-off between accuracy and computational efficiency. Results are shown in Table 4.

B. IMPACT OF FEATURE COMBINATION METHODS ON MULTI-CLASSIFICATION

When DenseNet201 features were used in conjunction with various ML models, we observed notable improvements. All are provided in Table 5 and Fig 3. KNN and CatBoost particularly excelled, achieving accuracies above 99%, with near-perfect precision and recall metrics. However, CatBoost required significantly more time (1246.11 seconds) compared to KNN's minimal processing time (0.02 seconds). XGBoost and LGBM also showed strong performance, both scoring high in accuracy and specificity. This underscores their

robustness in handling image classification tasks, making them suitable for applications where both accuracy and interpretability are important. DT, with its simpler structure, provided a solid accuracy of 93.36% and a specificity of 95%. Albeit not as high as the more complex models, its interpretability and ease of use make it valuable for scenarios where understanding the decision-making process is important. RF showed improved performance over individual DT, achieving an accuracy of 97.05% and specificity of 98%. This enhancement is attributed to RF's ability to reduce overfitting, a common issue in single DT, thus improving the robustness and generalizability of the model [43]. The MultinomialNB and SVM models, albeit less effective than the aforementioned models, still contributed valuable insights. MultinomialNB, renowned for its simplicity and efficiency in handling high-dimensional data, achieved an accuracy of 92.11%. This suggests its potential utility in scenarios where computational resources are limited. SVM, renowned for its effectiveness in high-dimensional spaces, demonstrated a commendable accuracy of 97.8%. Its ability to find the optimal hyperplane for class separation makes it a strong candidate for complex classification tasks, albeit with a higher computational cost as indicated by its processing time.

The integration of contour features with DenseNet201 resulted in consistent performance across the models, with KNN again demonstrating exceptional accuracy and efficiency, in Table 6 and Fig 4. However, SVM's performance strongly declined in this scenario, dropping to an accuracy of 33.1% and an unspecified specificity, indicating a poor fit for this combination of features. Other models like LGBM, CatBoost, and XGBoost maintained their high performance, showing less sensitivity to the addition of contour features compared to SVM. DT and RF still showed respectable results. DT's performance was moderately effective, indicating its potential limitations in handling the added complexity of contour features. On the other hand, RF showed an improved accuracy over DT, benefiting from its ensemble nature. MultinomialNB faced challenges with the added complexity of the combined features, resulting in less competitive performance compared to other models. However, its computational efficiency remained an advantage.

Incorporating histogram features with DenseNet201 significantly boosted performance, as demonstrated in Table 7 and Fig 5. KNN achieved an impressive accuracy of 99.68% and a perfect specificity score. This combination also enhanced the performance of the DT and RF models, increasing their accuracy and specificity compared to

TABLE 5. Performance metrics results of DenseNet201 features in conjunction with ML models. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.9863333	0.99	104.07	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
CatBoost	0.9881666	0.99	1246.11	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
XGBoost	0.9866666	0.99	84.26	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
KNN	0.9901666	0.99	0.01	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
DT	0.9336666	0.95	30.61	0.93 MA	0.93 MA	0.93 MA
				0.93 WA	0.93 WA	0.93 WA
RF	0.9705	0.98	48.75	0.97 MA	0.97 MA	0.97 MA
				0.97 WA	0.97 WA	0.97 WA
MultiNB	0.9211666	0.92	0.08	0.92 MA	0.92 MA	0.92 MA
				0.92 WA	0.92 WA	0.92 WA
SVM	0.978	0.99	13.51	0.98 MA	0.98 MA	0.98 MA
				0.98 WA	0.98 WA	0.98 WA

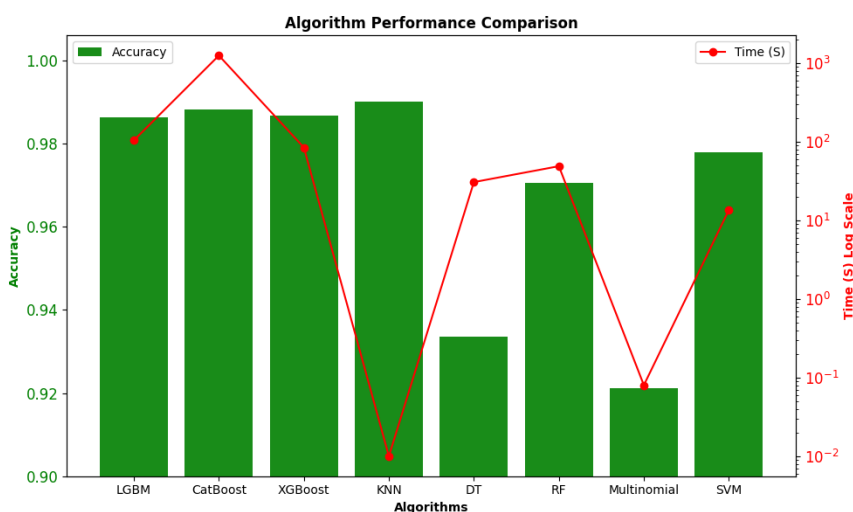


FIGURE 3. Accuracy and time results of DenseNet201 features in conjunction with ML models.

TABLE 6. Performance metrics results of DenseNet201 features integrated with contour features in conjunction with ML models. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.9863333	0.99	99.77	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
CatBoost	0.9868333	0.99	1253.89	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
XGBoost	0.9866666	0.99	83.17	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
KNN	0.9901666	0.99	0.01	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
DT	0.9315	0.95	31.95	0.93 MA	0.93 MA	0.93 MA
				0.93 WA	0.93 WA	0.93 WA
RF	0.970333	0.98	42.41	0.97 MA	0.97 MA	0.97 MA
				0.97 WA	0.97 WA	0.97 WA
MultiNB	0.9198333	0.92	0.18	0.92 MA	0.92 MA	0.92 MA
				0.92 WA	0.92 WA	0.92 WA
SVM	0.331	N/A	789.94	0.11 MA	0.33 MA	0.16 MA
				0.10 WA	0.32 WA	0.16 WA

using DenseNet201 features alone. The DT model notably improved its accuracy and specificity, indicating that the additional color textural information from histogram features complements its decision-making process. Similarly, the RF

capitalized on this combination, demonstrating enhanced accuracy and robustness compared to the DT. LGBM, XGBoost, and CatBoost also showed remarkable results with the added histogram features, utilizing the additional

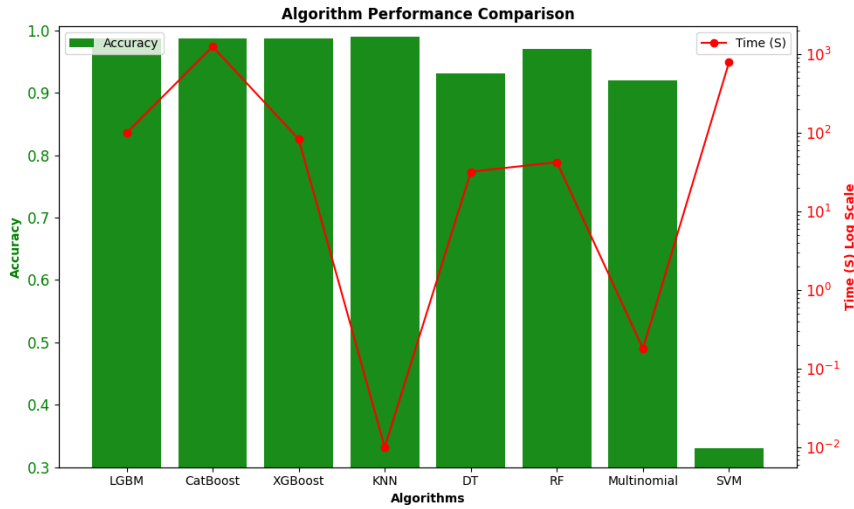


FIGURE 4. Accuracy and time results of DenseNet201 features integrated with contour features in conjunction with ML models.

TABLE 7. Performance metrics results of DenseNet201 features integrated with histogram features in conjunction with ML Models. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.9938333	0.99	117.84	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
CatBoost	0.9951666	1.00	1331.43	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
XGBoost	0.9948333	0.99	81.79	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
KNN	0.9968333	1.00	0.03	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
DT	0.946	0.96	30.61	0.95 MA	0.95 MA	0.95 MA
				0.95 WA	0.95 WA	0.95 WA
RF	0.9816666	0.98	33.26	0.98 MA	0.98 MA	0.98 MA
				0.98 WA	0.98 WA	0.98 WA
MultiNB	0.9243333	0.92	0.10	0.92 MA	0.92 MA	0.92 MA
				0.92 WA	0.92 WA	0.92 WA
SVM	0.9781666	0.99	80.44	0.98 MA	0.98 MA	0.98 MA
				0.98 WA	0.98 WA	0.98 WA

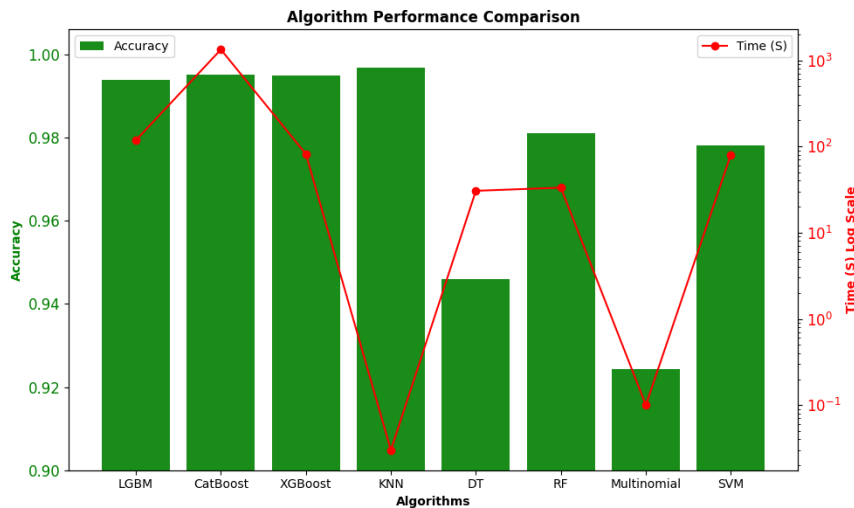


FIGURE 5. Accuracy and time results of DenseNet201 features integrated with histogram features in conjunction with ML models.

color textural information to improve their classification accuracy, demonstrating their adaptability and efficiency in

feature-rich environments. However, MultinomialNB consistently underperformed across all feature sets, suggesting

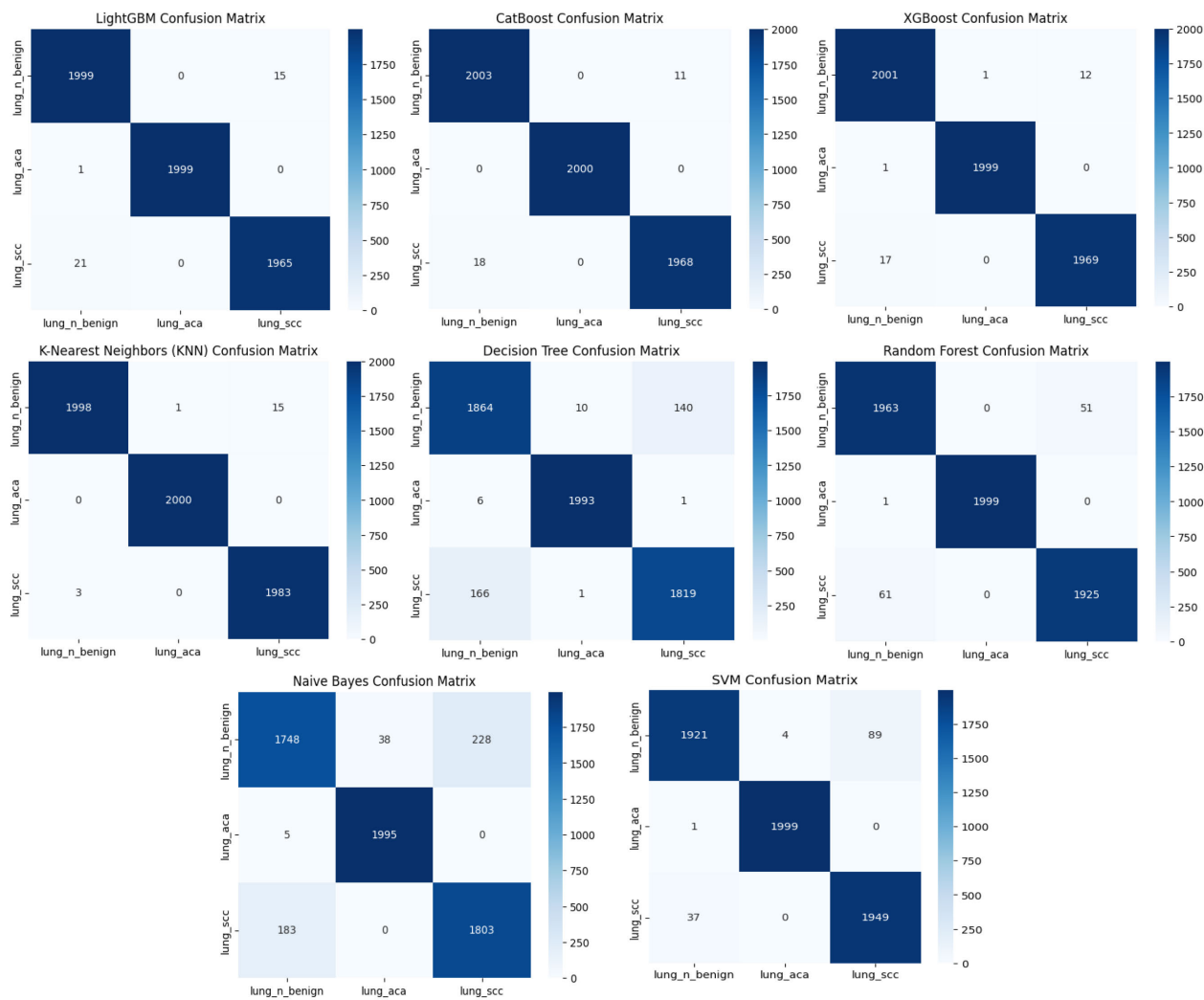


FIGURE 6. Confusion matrix of KNN by using the integrated features of DenseNet201 with histogram features.

its limitations for complex image classification tasks in this context. SVM showed respectable results with the added histogram features, indicating its ability to efficiently separate data in high-dimensional spaces, and suggesting its potential utility in scenarios where precision is critical, despite the computational demands.

The confusion matrices in Fig 6 display the performance of all ML models, which combine the DenseNet201 features with histogram features. To assess the model’s prediction power, these matrices show their instance classification performance across classes. The confusion matrix analysis reveals model strengths and weaknesses, misclassification trends, and classification approach improvements. Each confusion matrix provides insights into the true positives, false positives, true negatives, and false negatives for each class.

Notably, KNN misclassifies very few instances, with only 3 benign lung tissues misclassified as lung squamous cell carcinomas and almost perfect classification for lung

adenocarcinomas and lung squamous cell carcinomas. CatBoost and XGBoost models are effective in classifying lung adenocarcinomas without any error. LGBM performs well but struggles slightly with distinguishing between benign lung tissues and lung squamous cell carcinomas, as indicated by 15 misclassifications. Decision Tree and Naive Bayes models show significant misclassifications, especially between benign lung tissues and lung squamous cell carcinomas, and benign lung tissues and lung adenocarcinomas, respectively. This suggests these models have difficulty capturing the nuances between these classes. RF and SVM provide notable misclassifications in some areas, such as benign lung tissues being misclassified as lung squamous cell carcinomas.

The combined performance of these models, when used in conjunction with the features of DenseNet201, underscores the significance of integrating various types of FE methods and selecting the algorithm that is the most appropriate based on the specific requirements of the task. These requirements include accuracy, specificity, computational

TABLE 8. Performance metrics results of DenseNet201 features integrated with histogram features for binary classification lung benign vs. lung adenocarcinoma. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.999	1.00	32.34	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
CatBoost	0.9997	1.00	526.74	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
XGBoost	0.999	1.00	9.74	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
KNN	0.9987	1.00	0.01	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
DT	0.994	0.99	6.52	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
RF	0.999	1.00	12.79	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
MultiNB	0.99175	0.99	0.07	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
SVM	0.9985	0.99	2.76	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA

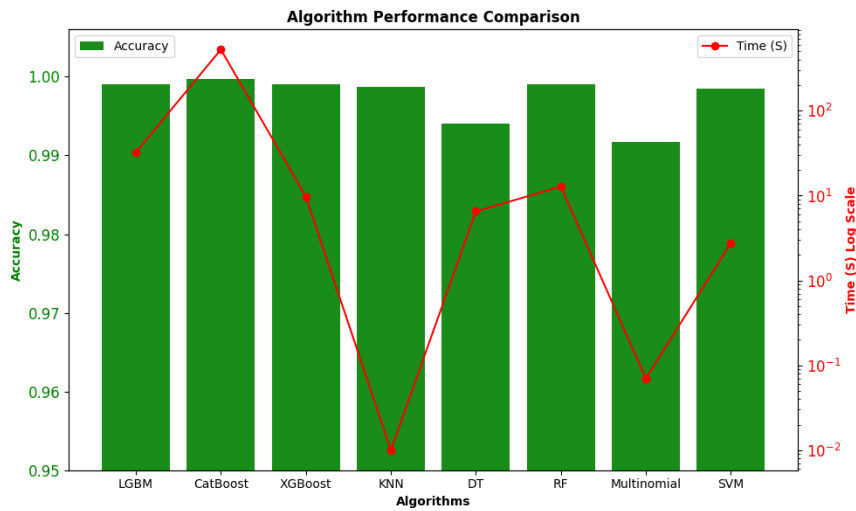


FIGURE 7. Accuracy and time of (DenseNet201 + histogram)features for binary classification lung benign vs. lung adenocarcinoma.

efficiency, and interpretability. The diversity in the capabilities of the algorithms provides a solid framework for solving the problems of LC image classification. Along with this, it offers several solutions to fit a wide range of clinical and research requirements. The results additionally underscore the potential of ML to improve the early detection of LC. The high accuracy and specificity achieved by our models, especially when utilizing DenseNet201 combined with histogram features, suggest that AI-driven approaches can significantly enhance the precision of LC diagnostics.

C. DENSENET201 FEATURES WITH HISTOGRAM FEATURES IMPACT ON BINARY CLASSIFICATION

In addition to multi-class classification, we conducted binary classification tasks to differentiate between specific LC subtypes and benign lung tissues. This aimed to refine our understanding of these subtypes’ unique characteristics.

1) LUNG BENIGN VS. LUNG ADENOCARCINOMA

In this classification task, most models achieved near-perfect performance metrics. Notably, KNN, LGBM, and CatBoost showed exceptional accuracy, specificity, and F1-scores, all reaching or exceeding 99.8%. CatBoost demonstrated the highest accuracy of 99.97% but required longer computation time (526.74 seconds). In contrast, KNN maintained its efficiency with an execution time of only 0.01 seconds. Results are demonstrated in Table 8 and Fig 7.

2) LUNG BENIGN VS. LUNG SQUAMOUS CELL CARCINOMA

In Table 9 and Fig 8, the results were even more remarkable in this classification scenario. All models, including KNN, LGBM, CatBoost, XGBoost, DT, RF, MultinomialNB, and SVM, achieved perfect accuracy and specificity scores of 100%. This demonstrates the unique ability of the models, especially when combined with DenseNet201 with histogram features, to distinguish benign lung tissue from squamous cell carcinoma.

TABLE 9. Performance metrics results of DenseNet201 features integrated with histogram features for binary classification lung benign vs. lung squamous cell carcinoma. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.999	1.00	47.25	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
CatBoost	1.0	1.00	528.88	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
XGBoost	1.0	1.00	7.69	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
KNN	1.0	1.00	0.01	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
DT	1.0	1.00	2.93	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
RF	1.0	1.00	6.72	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
MultiNB	1.0	1.00	0.07	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA
SVM	1.0	1.00	1.38	1.00 MA	1.00 MA	1.00 MA
				1.00 WA	1.00 WA	1.00 WA

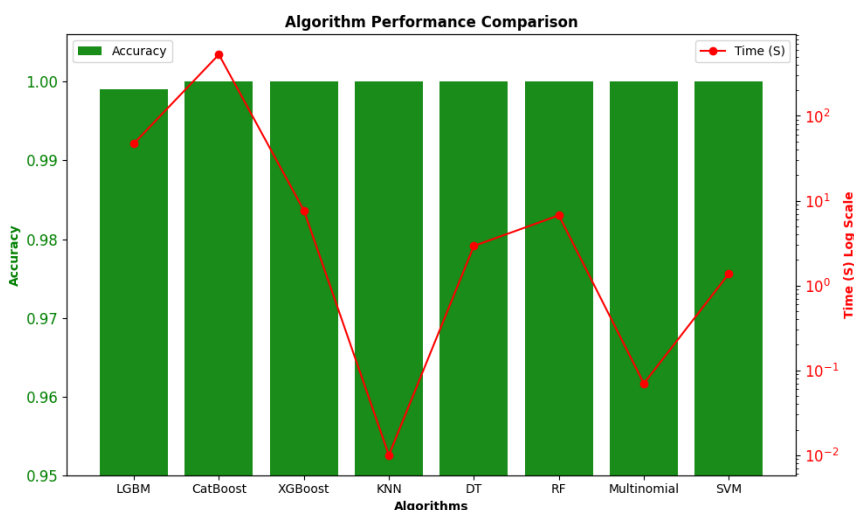


FIGURE 8. Accuracy and time of (DenseNet201 + histogram) features for binary classification lung benign vs. lung squamous cell carcinoma.

TABLE 10. Performance metrics results of DenseNet201 features integrated with histogram features for binary classification lung adenocarcinoma vs. lung squamous cell carcinoma. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Avg Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.988	0.99	53.08	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
CatBoost	0.99325	0.99	598.76	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
XGBoost	0.9895	0.99	38.52	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
KNN	0.98625	0.99	0.02	0.99 MA	0.99 MA	0.99 MA
				0.99 WA	0.99 WA	0.99 WA
DT	0.9213	0.92	31.93	0.92 MA	0.92 MA	0.92 MA
				0.92 WA	0.92 WA	0.92 WA
RF	0.9755	0.98	27.11	0.98 MA	0.98 MA	0.98 MA
				0.98 WA	0.98 WA	0.98 WA
MultiNB	0.89325	0.89	0.12	0.89 MA	0.89 MA	0.89 MA
				0.89 WA	0.89 WA	0.89 WA
SVM	0.975	0.97	15.07	0.97 MA	0.98 MA	0.97 MA
				0.98 WA	0.97 WA	0.98 WA

3) LUNG ADENOCARCINOMA VS. LUNG SQUAMOUS CELL CARCINOMA

From Table 10 and Fig 9, differentiating between adenocarcinoma and squamous cell carcinoma posed a slightly greater

challenge. However, the models still performed admirably, with CatBoost leading at 99.32% accuracy. KNN, LGBM, and XGBoost also showcased high accuracy above 98.6%. DTs and MultinomialNB, albeit less accurate than other

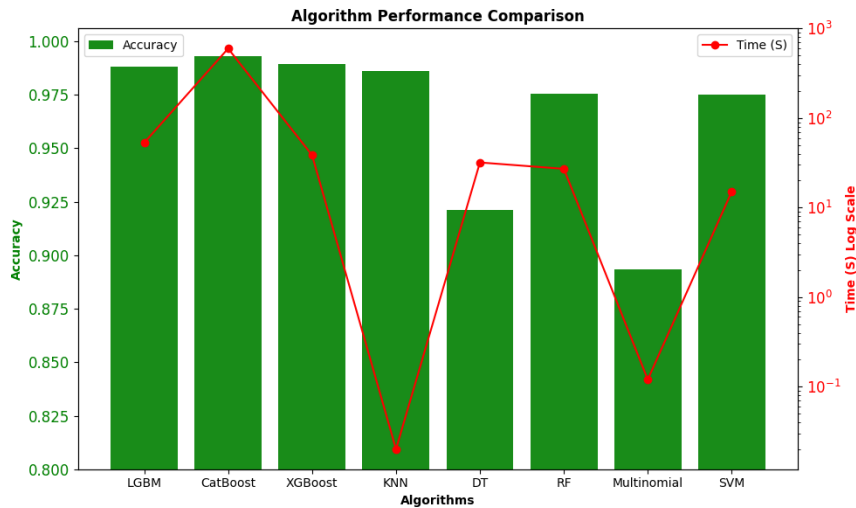


FIGURE 9. Accuracy and time of (DenseNet201 + histogram) features for binary classification lung adenocarcinoma vs. lung squamous cell carcinoma.

models, still offered valuable insights with accuracies above 89%.

The binary classification results emphasize the nuanced capabilities of these models in discerning between specific LC types and benign conditions. The consistently high performance across various models, especially in the benign vs. adenocarcinoma and benign vs. squamous cell carcinoma tasks, demonstrates the potential of ML algorithms in highly specific and accurate LC diagnosis. The relative decrease in performance in distinguishing between adenocarcinoma and squamous cell carcinoma underscores the complexity of this task. Nevertheless, the models, particularly CatBoost, managed to maintain high accuracy levels, reinforcing the efficacy of the chosen FE methods. All that marks a significant step towards AI-assisted diagnostics in oncology.

D. DATA VISUALIZATION

The scatter plot in Fig 10, shows the dataset projected onto the principal components Analysis (PCA). PCA reduces the dimensionality of the data by finding the axes along which the variance is maximized. Remember, reducing dimensions might oversimplify the dataset, and important features might be lost in this reduction, which may not capture the full complexity of the data. In this plot, the 3D PCA represent the directions in the dataset that account for the most variance. The distribution of points gives insights into the separability of the data. This can partly explain why certain algorithms, like KNN, performed well if they are effectively capturing these separable structures in higher-dimensional space. The first three principal components together capture about 34.3% of the total variance in the data, with the first component accounts for approximately 21.3%, the second component for about 7.1%, and the third component adds another 5.9%. This indicates that while these components reveal some structure in the data, a significant portion of the variance (over 65.7%) is not captured in this three-dimensional representation.

For a detailed analysis of cluster formation and overlap, we typically rely on visual inspection. The moderate level of variance captured suggests that while PCA provides some insights, other components (not visualized here) might also contain important information for classification. The clusters appear distinct, it may indicate that the categories are well-separated in the feature space, potentially explaining the good performance of KNN, which performs well when similar instances are closer together in the feature space.

E. METHODOLOGY VALIDATION ACROSS DIFFERENT DATASET

In order to assess the versatility and generalizability of our ML algorithms, we extended our analysis to another crucial area in oncology: breast cancer. Utilizing the BreakHis dataset, which consists small number of histopathological images, we explored the performance of the same set of algorithms that were applied to the LC dataset. This cross-application aims to understand how well the methodologies and insights gleaned from LC classification can be transferred to another context within oncology. Utilizing the same DenseNet201 + Histogram FE method as used in our LC data. We conducted binary classifications to distinguish between malignant and benign breast cancer cases. Below are the summarized results of each algorithm's performance on this dataset: KNN demonstrated high accuracy and specificity, with minimal processing time, reinforcing its efficiency and effectiveness in image classification tasks across different cancer types. LGBM matched KNN in accuracy and specificity, albeit with a longer processing time, underscoring its robustness in handling complex image data. XGBoost equaled KNN and LGBM in accuracy and specificity, indicating its strong adaptability and efficiency in diverse medical imaging contexts. CatBoost showed slightly lower accuracy compared to KNN, LGBM, and XGBoost but with substantially higher computational demands, suggesting

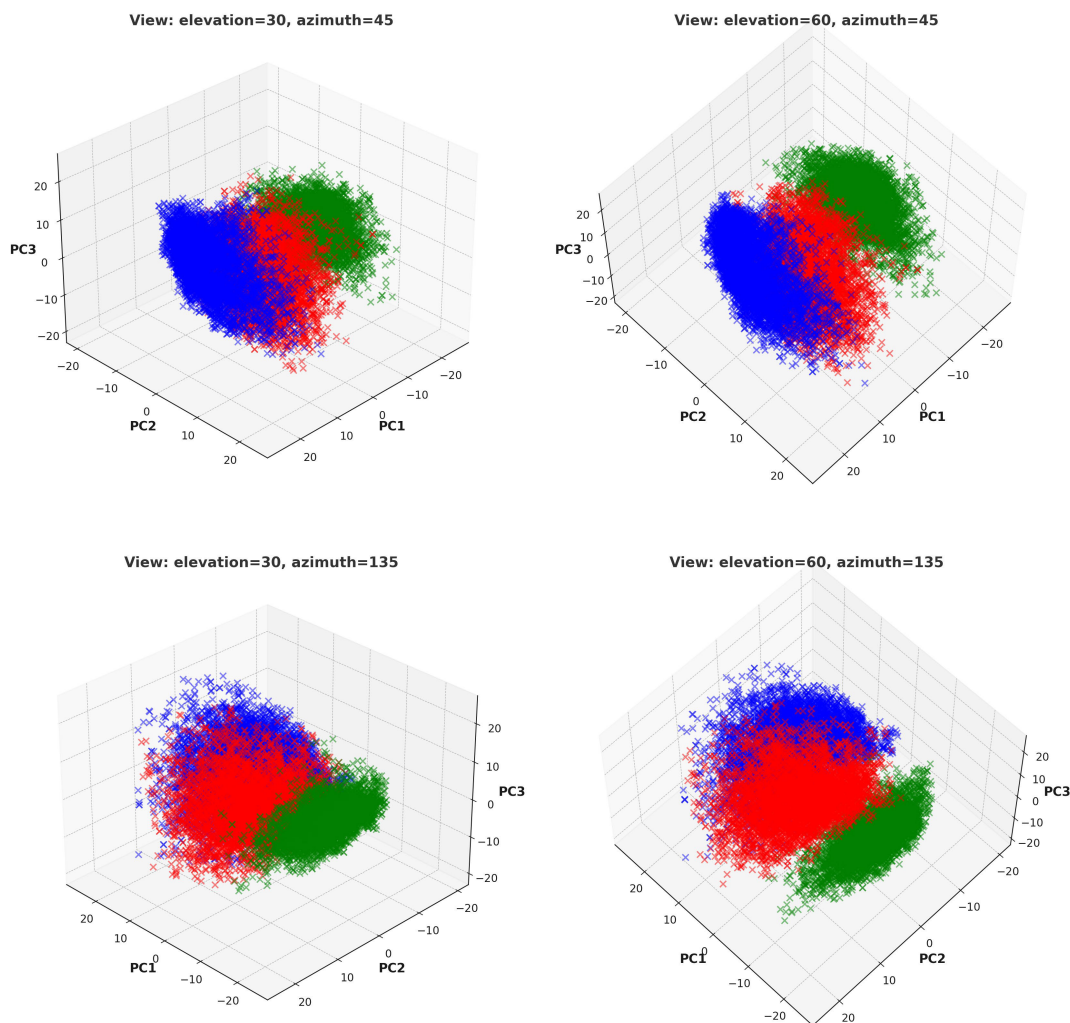


FIGURE 10. Multiple perspectives of 3D PCA data visualization scatter plot obtained from the extracted features of the model.

TABLE 11. The performance metrics results of DenseNet201 features integrated with histogram features for binary classification breast cancer. MA (micro average), WA (weighted average).

Algorithm	Accuracy	Specificity	Time (S)	Precision	Recall	F1-score
LGBM	0.948087	0.94	35.73	0.94 MA	0.94 MA	0.94 MA
				0.95 WA	0.95 WA	0.95 WA
CatBoost	0.939890	0.94	667.67	0.94 MA	0.92 MA	0.93 MA
				0.94 WA	0.94 WA	0.94 WA
XGBoost	0.948087	0.94	42.62	0.94 MA	0.94 MA	0.94 MA
				0.95 WA	0.95 WA	0.95 WA
KNN	0.948087	0.96	0.02	0.96 MA	0.93 MA	0.94 MA
				0.95 WA	0.95 WA	0.95 WA
DT	0.871584	0.86	5.59	0.86 MA	0.85 MA	0.85 MA
				0.87 WA	0.87 WA	0.87 WA
RF	0.928961	0.92	6.44	0.92 MA	0.92 MA	0.92 MA
				0.93 WA	0.93 WA	0.93 WA
MultiNB	0.860655	0.84	0.05	0.84 MA	0.84 MA	0.84 MA
				0.86 WA	0.86 WA	0.86 WA
SVM	0.931693	0.85	17.16	0.93 MA	0.91 MA	0.92 MA
				0.93 WA	0.93 WA	0.93 WA

a trade-off between performance and efficiency. DT, while less accurate and specific than other models, its faster execution time and interpretability remain advantageous,

particularly in scenarios where rapid results are essential. RF Exhibited improved accuracy and specificity over DT. MultinomialNB although had the lowest accuracy and

TABLE 12. Comparative result of the proposed method with other related works.

Work	Year	Accuracy	Model
[13]	2022	98.60%	DenseNet121 FE + RF
[20]	2022	98.4%	AlexNet CNN+ Histogram Equalization
[21]	2020	97.92%	shallow CNN
[22]	2022	97.11%	Customized CNN
[23]	2022	94.42%	XGBoost
[24]	2023	89%	Multi-level CNN
[25]	2023	91.57%	GWO FE + IWO Feature Selection + hyperparameter tuning RAdam + DT
[26]	2023	98.9%	HOG FE + hyperparameter tuning GAO + IGNN
Proposed Model	2024	99.68333%	DenseNet201 + Color Histogram + KNN

specificity, its extremely quick execution time positions it as a viable option in resource-constrained settings. SVM presented respectable accuracy and specificity, with moderate computational demands, highlighting its potential in high-dimensional data scenarios where precision is paramount.

The findings on the various performance measures in clinical and research contexts, in particular, demonstrate how adaptable and versatile our chosen ML classifiers are when it comes to dealing with the different types of cancer. KNN, LGBM, and XGBoost have all demonstrated consistently good performance across both datasets, which shows that these algorithms have great promise in a wide range of histopathological image classification tasks.

IV. DISCUSSION

The study discusses the use of advanced AI methodologies fusion for the analysis of histopathological images in the context of early detection and classification of LC. The significant findings offer a promising direction due to the high mortality rate and challenges in the early diagnosis of LC. The study highlights the performance of ML models, specifically KNN, CatBoost, XGBoost, and LGBM, in multi-class and binary classification tasks. It emphasizes the importance of algorithm selection based on requirements such as accuracy, specificity, and computational efficiency. For instance, KNN demonstrates high accuracy and efficiency in both multi-class and binary classification tasks with minimal misclassification, making it suitable for real-time applications prioritizing accuracy and speed. On the other hand, CatBoost and XGBoost exhibit high accuracy despite longer computation times, indicating its potential for scenarios where precision is paramount. For applications where prediction time can be longer, and the focus is on minimizing false positives, CatBoost or XGBoost may be preferred. Decision Tree, Naive Bayes, Random Forest, and SVM models show lower performance compared to KNN, CatBoost, and XGBoost, especially in terms of misclassification rates and, for some, in computational efficiency. The study demonstrates the value of integrating DenseNet201 with contour and histogram features, highlighting the benefits of combining different types of FE methods to enhance the power of classifiers. This hybrid approach shows promise in advancing the precision of LC diagnostics. The findings also reveal a trade-off between computational efficiency and classification accuracy. The

use of models like CatBoost, while offering high accuracy, comes with the trade-off of longer computational times, emphasizing the need for optimization strategies in real-world applications where both accuracy and efficiency are crucial. Additionally, the challenges in distinguishing between adenocarcinoma and squamous cell carcinoma in LC subtypes reveal the complexity inherent in their classification, highlighting the necessity for further research to refine algorithms and FE methods for more nuanced differentiation.

The study compares the proposed model with other state-of-the-art studies in classifying LC histopathological images, as shown in Table 12. The presented model stands out with an impressive accuracy of 99.68%, surpassing other models. This superior performance is credited to the effective combination of advanced FE using DenseNet201 and the robust classification capabilities of KNN. This combination not only boosts the model's accuracy but also ensures computational efficiency. Unlike other methods, such as AlexNet CNN with Histogram Equalization or DenseNet121 FE with RF, the model shows an improvement in accuracy, showcasing the effectiveness of our hybrid FE method. This study introduces a novel approach by integrating DenseNet201 with color histogram features, a method not explored in the referenced studies. This innovation enhances the model's ability to detect subtle variations in histopathological images, which is crucial for the early detection of LC.

A. RESEARCH GAP, LIMITATIONS, AND FUTURE WORK

This research tackles a significant gap in the current scholarly discussion on the precision and speed of lung cancer (LC) diagnosis through histopathological images. Despite progress in diagnostic technologies, the issue of quickly and precisely identifying lung adenocarcinomas, lung squamous cell carcinomas, and non-cancerous lung tissues with minimal manual oversight remains unresolved. The study employs DenseNet201, color histogram methods, and various machine learning (ML) techniques to improve diagnostic capabilities, addressing an important gap in automated image-based cancer detection.

Nonetheless, the methodology presented is not flawless. The LC25000 dataset, while extensive, might not capture the full range of variability seen in wider clinical environments. Depending on this dataset could raise questions about the applicability of the results in various imaging scenarios. This

reliance poses a risk of selection bias, given the dataset may not adequately reflect the range of LC histopathologies seen in medical practice. The performance metrics of the algorithms, are particular to the datasets utilized, and the extension of these findings to different cancers or their subtypes requires further verification.

To overcome these shortcomings, future research should aim to diversify the dataset to encompass a wider array of histopathological images. This would help reduce selection bias and strengthen the model's robustness. Moreover, investigating different data augmentation methods could decrease processing time, enhancing feasibility for real-time clinical use. Implementing different image preprocessing can lessen the influence of image quality on the outcomes. Employing feature selection techniques to focus on the most relevant features could refine the model, boosting both its effectiveness and efficiency. Utilizing image segmentation to more accurately target areas of interest in the images could improve the model's ability to identify subtle indicators of cancer. Investigating the use of advanced FE methods might lead to a greater sensitivity in detecting complex patterns in histopathological images. These improvements are crucial for increasing the diagnostic precision and clinical utility of ML models in LC diagnosis, ensuring they fulfill the stringent requirements of medical practice.

V. CONCLUSION

The presented research aimed to improve the early detection and classification of lung cancer (LC) by using advanced AI methodologies. It combined DenseNet201 for deep feature extraction (FE) with color histogram features and analyzed them using various machine learning (ML) algorithms, particularly KNN, which showed exceptional performance. Our model, tested on the LC25000 dataset, achieved a remarkable accuracy of 99.68%, outperforming state-of-the-art models. This high accuracy is crucial due to the high mortality rate and the challenges in the early diagnosis of LC. The study also highlights the importance of selecting the right algorithm for specific needs, such as KNN for real-time applications due to its efficiency and accuracy and CatBoost for accurate applications despite longer computation periods. Furthermore, the integration of multiple FE approaches not only improved classifier discrimination but also showed adaptability to other cancers, as demonstrated by its application to the BreakHis dataset of histopathological images for breast cancer. This suggests its potential for use in other critical areas of oncology. It clearly recognizes the trade-offs between computational efficiency and classification accuracy and the difficulties of classifying LC subtypes. This emphasizes the need to optimize and develop algorithms and FE approaches.

Our paper introduces LC diagnosis in a novel way, promising the future of AI-driven oncology diagnostics. It opens new research and development options, especially for more delicate and effective cancer diagnosis and classification.

As we grow in this discipline, data scientists, medical specialists, and oncologists must collaborate to translate technical advances into therapeutic applications.

CONFLICTS OF INTEREST

The authors declare no conflict of interest.

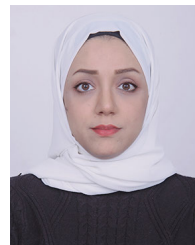
DECLARATION OF COMPETING INTEREST

The authors declare that this manuscript is original, has not been published before, and is not currently being considered for publication elsewhere.

REFERENCES

- [1] R. L. Siegel, K. D. Miller, N. S. Wagle, and A. Jemal, "Cancer statistics, 2023," *CA, Cancer J. Clinicians*, vol. 73, no. 1, pp. 17–48, Jan. 2023, doi: [10.3322/caac.21763](https://doi.org/10.3322/caac.21763).
- [2] Amer. Cancer Soc. (Oct. 18, 2023). *Cancer Statistics Center, American Cancer Society*. Accessed: Dec. 5, 2023. [Online]. Available: <https://cancerstatisticscenter.cancer.org/#/>
- [3] F. Venuta, D. Diso, I. Onorati, M. Anile, S. Mantovani, and E. A. Rendina, "Lung cancer in elderly patients," *J. Thoracic Disease*, vol. 8, no. S11, pp. S908–S914, Nov. 2016, doi: [10.21037/jtd.2016.05.20](https://doi.org/10.21037/jtd.2016.05.20).
- [4] A. Leiter, R. R. Veluswamy, and J. P. Wisnivesky, "The global burden of lung cancer: Current status and future trends," *Nature Rev. Clin. Oncol.*, vol. 20, no. 9, pp. 624–639, Sep. 2023, doi: [10.1038/s41571-023-00798-3](https://doi.org/10.1038/s41571-023-00798-3).
- [5] (Jun. 26, 2023). *World Health Organization: WHO*. Accessed: Jan. 4, 2024. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/lung-cancer>
- [6] Lung Cancer Statistics | CDC. *Centers for Disease Control and Prevention*. Accessed: Feb. 4, 2024. [Online]. Available: <https://www.cdc.gov/cancer/lung/statistics/index.htm>
- [7] SEER. (2023). *Common Cancer Sites—Cancer Stat Facts*. Accessed: Dec. 20, 2023. [Online]. Available: <https://seer.cancer.gov/statfacts/html/common.html>
- [8] A. J. Tybjerg, S. Friis, K. Brown, M. C. Nilbert, L. Mørch, and B. Køster, "Updated fraction of cancer attributable to lifestyle and environmental factors in Denmark in 2018," *Sci. Rep.*, vol. 12, no. 1, p. 549, Jan. 2022, doi: [10.1038/s41598-021-04564-2](https://doi.org/10.1038/s41598-021-04564-2).
- [9] M. Del Re, E. Rofi, G. Restante, S. Crucitta, E. Arrigoni, S. Fogli, M. Di Maio, I. Petrini, and R. Danesi, "Implications of Kras mutations in acquired resistance to treatment in NSCLC," *Oncotarget*, vol. 9, no. 5, pp. 6630–6643, Jan. 2018, doi: [10.18632/oncotarget.23553](https://doi.org/10.18632/oncotarget.23553).
- [10] (Mar. 22, 2022). *Lung Cancer—Diagnosis and Treatment, Mayo Clinic*. Accessed: Jan. 11, 2024. [Online]. Available: <https://www.mayoclinic.org/diseases-conditions/lung-cancer/diagnosis-treatment/drc-20374627>
- [11] M. Al-Jabbar, M. Alshahrani, E. M. Senan, and I. A. Ahmed, "Histopathological analysis for detecting lung and colon cancer malignancies using hybrid systems with fused features," *Bioengineering*, vol. 10, no. 3, p. 383, Mar. 2023, doi: [10.3390/bioengineering10030383](https://doi.org/10.3390/bioengineering10030383).
- [12] S. Garg and S. Garg, "Prediction of lung and colon cancer through analysis of histopathological images by utilizing pre-trained CNN models with visualization of class activation and saliency maps," in *Proc. 3rd Artif. Intell. Cloud Comput. Conf.*, New York, NY, USA, Dec. 2020, pp. 38–45, doi: [10.1145/3442536.3442543](https://doi.org/10.1145/3442536.3442543).
- [13] N. Kumar, M. Sharma, V. P. Singh, C. Madan, and S. Mehandia, "An empirical study of handcrafted and dense feature extraction techniques for lung and colon cancer classification from histopathological images," *Biomed. Signal Process. Control*, vol. 75, May 2022, Art. no. 103596, doi: [10.1016/j.bspc.2022.103596](https://doi.org/10.1016/j.bspc.2022.103596).
- [14] (Jun. 19, 2019). *How Do Cancer Cells Grow and Spread?* Accessed: Dec. 5, 2023. [Online]. Available: <https://ncbi.nlm.nih.gov/books/NBK279410/>
- [15] R. W. Pettit, J. Byun, Y. Han, Q. T. Ostrom, J. Edelson, K. M. Walsh, M. L. Bondy, R. J. Hung, J. D. McKay, and C. I. Amos, "The shared genetic architecture between epidemiological and behavioral traits with lung cancer," *Sci. Rep.*, vol. 11, no. 1, p. 17559, Sep. 2021, doi: [10.1038/s41598-021-96685-x](https://doi.org/10.1038/s41598-021-96685-x).

- [16] M. Masud, N. Sikder, A.-A. Nahid, A. K. Bairagi, and M. A. AlZain, "A machine learning approach to diagnosing lung and colon cancer using a deep learning-based classification framework," *Sensors*, vol. 21, no. 3, p. 748, Jan. 2021, doi: [10.3390/s21030748](https://doi.org/10.3390/s21030748).
- [17] M. Ali and R. Ali, "Multi-input dual-stream capsule network for improved lung and colon cancer classification," *Diagnostics*, vol. 11, no. 8, p. 1485, Aug. 2021, doi: [10.3390/diagnostics11081485](https://doi.org/10.3390/diagnostics11081485).
- [18] N. Baranwal, P. Doravari, and R. Kachhoria, "Classification of histopathology images of lung cancer using convolutional neural network (CNN)," in *Disruptive Developments in Biomedical Applications (CNN)*. Boca Raton, FL, USA: CRC Press, 2022, pp. 75–89.
- [19] B. K. Hatuwal and H. C. Thapa, "Lung cancer detection using convolutional neural network on histopathological images," *Int. J. Comput. Trends Technol.*, vol. 68, no. 10, pp. 21–24, Oct. 2020, doi: [10.14445/22312803/ijctt-v68i10p104](https://doi.org/10.14445/22312803/ijctt-v68i10p104).
- [20] S. Mehmood, T. M. Ghazal, M. A. Khan, M. Zubair, M. T. Naseem, T. Faiz, and M. Ahmad, "Malignancy detection in lung and colon histopathology images using transfer learning with class selective image processing," *IEEE Access*, vol. 10, pp. 25657–25668, 2022, doi: [10.1109/ACCESS.2022.3150924](https://doi.org/10.1109/ACCESS.2022.3150924).
- [21] S. Mangal, A. Chaurasia, and A. Khajanchi, "Convolution neural networks for diagnosing colon and lung cancer histopathological images," 2020, *arXiv:2009.03878*.
- [22] J. Civit-Masot, A. Bañuls-Beaterio, M. Domínguez-Morales, M. Rivas-Pérez, L. Muñoz-Saavedra, and J. M. Rodríguez Corral, "Non-small cell lung cancer diagnosis aid with histopathological images using explainable deep learning techniques," *Comput. Methods Programs Biomed.*, vol. 226, Nov. 2022, Art. no. 107108, doi: [10.1016/j.cmpb.2022.107108](https://doi.org/10.1016/j.cmpb.2022.107108).
- [23] M. Mamun, A. Farjana, M. Al Mamun, and M. S. Ahammed, "Lung cancer prediction model using ensemble learning techniques and a systematic review analysis," in *Proc. IEEE World AI IoT Congr.*, Jul. 2022, pp. 187–193, doi: [10.1109/AIIoT54504.2022.9817326](https://doi.org/10.1109/AIIoT54504.2022.9817326).
- [24] M. Ramesh, S. Maheswaran, S. Theivanayaki, K. Kodeeswari, L. Krishnasamy, and N. Sriram, "Efficient lung cancer classification on multi level convolution neural network using histopathological images," in *Proc. 14th Int. Conf. Comput. Commun. Netw. Technol. (ICCCNT)*, Jul. 2023, pp. 1–7, doi: [10.1109/icccnt56998.2023.10307852](https://doi.org/10.1109/icccnt56998.2023.10307852).
- [25] K. Shanmugam and H. Rajaguru, "Exploration and enhancement of classifiers in the detection of lung cancer from histopathological images," *Diagnostics*, vol. 13, no. 20, p. 3289, Oct. 2023, doi: [10.3390/diagnostics13203289](https://doi.org/10.3390/diagnostics13203289).
- [26] S. Dinesh Krishnan, D. Pelusi, A. Daniel, V. Suresh, and B. Balusamy, "Improved graph neural network-based green anaconda optimization for segmenting and classifying the lung cancer," *Math. Biosci. Eng.*, vol. 20, no. 9, pp. 17138–17157, Sep. 2023, doi: [10.3934/mbe.2023764](https://doi.org/10.3934/mbe.2023764).
- [27] S. Chakraborty and K. Mali, "An overview of biomedical image analysis from the deep learning perspective," in *Research Anthology on Improving Medical Imaging Techniques for Analysis and Intervention*. Hershey, PA, USA: IGI Global, 2023, ch. 3, pp. 43–59, doi: [10.4018/978-1-6684-7544-7.ch003](https://doi.org/10.4018/978-1-6684-7544-7.ch003).
- [28] A. A. Borkowski, M. M. Bui, L. Brannon Thomas, C. P. Wilson, L. A. DeLand, and S. M. Mastorides, "Lung and colon cancer histopathological image dataset (LC25000)," 2019, *arXiv:1912.12142*.
- [29] R. C. Gonzalez, *Digital Image Processing*, 3rd ed. Upper Saddle River, NJ, USA: Pearson, 2009.
- [30] D. Singh and B. Singh, "Investigating the impact of data normalization on classification performance," *Appl. Soft Comput.*, vol. 97, Dec. 2020, Art. no. 105524, doi: [10.1016/j.asoc.2019.105524](https://doi.org/10.1016/j.asoc.2019.105524).
- [31] A. Güneş, H. Kalkan, and E. Durmuş, "Optimizing the color-to-grayscale conversion for image classification," *Signal, Image Video Process.*, vol. 10, no. 5, pp. 853–860, Jul. 2016.
- [32] M. Sonka, V. Hlavac, and R. Boyle, *Image Processing, Analysis, and Machine Vision*. Stamford, CT, USA: Cengage Learning, 2014.
- [33] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 2261–2269, doi: [10.1109/CVPR.2017.243](https://doi.org/10.1109/CVPR.2017.243).
- [34] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep Learning*, vol. 1. Cambridge, MA, USA: MIT Press, 2016.
- [35] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning*, vol. 112. New York, NY, USA: Springer, 2013. [Online]. Available: <https://link.springer.com/content/pdf/10.1007/978-1-0716-1418-1.pdf>
- [36] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu, "LightGBM: A highly efficient gradient boosting decision tree," in *Advances in Neural Information Processing Systems*, vol. 30, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Red Hook, NY, USA: Curran Associates, 2017. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf
- [37] A. Veronika Dorogush, V. Ershov, and A. Gulin, "CatBoost: Gradient boosting with categorical features support," 2018, *arXiv:1810.11363*.
- [38] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, San Francisco, CA, USA, Aug. 2016, pp. 785–794, doi: [10.1145/2939672.2939785](https://doi.org/10.1145/2939672.2939785).
- [39] K.-A.-N. Nguyen, P. Tandon, S. Ghanavati, S. N. Cheetirala, P. Timsina, R. Freeman, D. Reich, M. A. Levin, M. Mazumdar, Z. A. Fayad, and A. Kia, "A hybrid decision tree and deep learning approach combining medical imaging and electronic medical records to predict intubation among hospitalized patients with COVID-19: Algorithm development and validation," *JMIR Formative Res.*, vol. 7, Oct. 2023, Art. no. e46905, doi: [10.2196/46905](https://doi.org/10.2196/46905).
- [40] L. Breiman, "Random Forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, Oct. 2001, doi: [10.1023/a:1010933404324](https://doi.org/10.1023/a:1010933404324).
- [41] A. McCallum and K. Nigam, "A comparison of event models for naive Bayes text classification," in *Proc. AAAI Workshop Learn. Text Categorization*, 1998, vol. 752, no. 1, pp. 41–48. [Online]. Available: <http://yangli-feasibility.com/home/classes/lfd2022fall/media/aaaiws98.pdf>
- [42] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, Sep. 1995, doi: [10.1007/bf00994018](https://doi.org/10.1007/bf00994018).
- [43] (Jan. 3, 2024). *Understand Random Forest Algorithms With Examples (Updated 2024)*. Accessed: Jan. 11, 2024. [Online]. Available: <https://www.analyticsvidhya.com/blog/2021/06/understanding-random-forest/>



NAGLAA F. NOAMAN was born in Taiz, Yemen, in February 1992. She received the bachelor's degree in software engineering from the Faculty of Engineering and Information Technology, Taiz University, Taiz, and the Master of Engineering degree in computer science and technology from Xidian University, Xi'an, China, where she is currently pursuing the Ph.D. degree with the School of Artificial Intelligence. She is a devoted professional within the realm of computer science and technology. Her technological odyssey commenced in Taiz. Her academic trajectory has been characterized by an unwavering commitment to excellence. Her research interest includes the study of lung cancer. She is wholeheartedly dedicated to pushing the boundaries of knowledge in this dynamic and vital field.



BASSAM M. KANBER was born in Taiz, Yemen, in March 1991. He received the bachelor's degree in software engineering from the Faculty of Engineering and Information Technology, Taiz University, Taiz, and the Master of Engineering degree in computer science and technology from Xidian University, Xi'an, China, where he is currently pursuing the Ph.D. degree with the School of Artificial Intelligence. His academic path has been marked by a commitment to excellence. His journey in the world of technology began in Taiz. He is a dedicated professional in the field of computer science and technology. His research is centered around breast cancer and he is dedicated to pushing the boundaries of knowledge in this dynamic field.



AHMAD AL SMADI received the Ph.D. degree in computer science and technology from the School of Artificial Intelligence, Xidian University, Xi’an, China. He was a Research Assistant with the College of Technological Innovation, Zayed University, United Arab Emirates. He is currently an Assistant Professor with the Department of Data Science and Artificial Intelligence, Zarqa University. His research interests include information systems, computer vision, machine learning, and deep learning.



MUTASEM K. ALSMADI (Member, IEEE) is currently an Associate Professor with the College of Applied Studies and Community Service. He has made significant contributions to the field of research, particularly in the areas of image segmentation and learning management. With an H-index of 36, he has achieved substantial recognition for his scholarly work. Throughout his academic career, he has coauthored 81 publications that have garnered a remarkable total of 4,101 citations. His research output demonstrates both breadth and depth in addressing key challenges and advancing knowledge within the fields of image segmentation and learning management.

...



LICHENG JIAO (Fellow, IEEE) received the B.S. degree from Shanghai Jiao Tong University, Shanghai, China, in 1982, and the M.S. and Ph.D. degrees from Xi’an Jiaotong University, Xi’an, China, in 1984 and 1990, respectively.

Since 1992, he has been a Distinguished Professor with the School of Electronic Engineering, Xidian University, Xi’an, where he is currently the Director of the Key Laboratory of Intelligent Perception and Image Understanding of the Ministry of Education of China. His research interests include machine learning, deep learning, natural computation, remote sensing, image processing, and intelligent information processing.

Prof. Jiao has been a Foreign Member of the Academia European and Russian Academy of Natural Sciences. He is the Chairman of the Awards and Recognition Committee, the Vice Board Chairperson of the Chinese Association of Artificial Intelligence, a fellow of IET/CAAI/CIE/CCF/CAA/CSIG/AAIA/AIIA/ACIS, the Councilor of Chinese Institute of Electronics, a Committee Member of Chinese Committee of Neural Networks, and an Expert of the Academic Degrees Committee of the State Council.